
CHAPTER 5

Further characterisation of the *intB* element in *D. nodosus* strain A198

5.1 Introduction

There are two copies, or part thereof, of the *intB* element in the genome of the virulent *D. nodosus* strain A198 (Bloomfield *et al.*, 1997b). One copy of the *intB* element is present adjacent to the *attR* site of *vap* region 3, and sequencing analyses to date have led to the identification of five potential open reading frames, designated *intB*, *regA*, *gepA*, *gepB* and *gepC* (Table 5.1) (Bloomfield *et al.*, 1997b). The location of this genetic element in relation to the *vap* regions is shown in Figure 5.1.

The second copy of the *intB* element sequences are present next to the *attR* site of *vap* region 2 in strain A198, and contains a partial copy of the *intB* gene, encoding of only the first 164 aa from the amino-terminus of the 403 aa *intB* gene, and is called *intB_N* (Figure 5.1) (Bloomfield *et al.*, 1997b; Whittle, 1994). The *intB_N* gene is followed by *orf379*, which is not related to *regA* or *gepB-C*.

In benign *D. nodosus* strain C305, there are two copies of the *intB* gene, both of which are truncated. One of these copies of the *intB* gene is integrated adjacent to *pnpA* in the C305 genome (Figure 5.1), and is identical to *intB_N* from *D. nodosus* strain A198 (Bloomfield, 1997; Shaw, 1997). This is followed by *orf379*. The other copy of *intB* is next to the *intC* element sequences and consists of only 100 bp from the middle of the *intB*

Table 5.1: Sequence analysis of the *intB* element^a in *D. nodosus* strain A198

Gene	Co-ordinates 5'-3- (nt)	Size (aa)	% identity to nt or putative protein	Homologues/Description	Accession Number	P(π)				
intB	389-1599	403	33.1/314aa	<i>Coli</i> retronphage phiR73 integrase [<i>selCIRNA</i>]	M64113	7e-31				
				<i>E. coli</i> 0157:H7 LEE pathogenicity island integrase [<i>IRNAsec</i>]	AF071034	6.9e-34				
				<i>Mesorhizobium loti</i> symbiosis island integrase [<i>IRNAphe</i>]	AF049242	2.1e-27				
				<i>Vibrio cholerae</i> 0395 integrase from CTXφ PAI [<i>ssrA</i>]	U02372	4.3e-27				
				<i>E. coli</i> K12 cryptic prophage P4-57 integrase [<i>ssrA</i>]	U11296	2.4e-26				
				<i>Shigella flexneri</i> φSF6 integrase	X59553	2.7e-26				
				<i>E. coli</i> K12 integrase [<i>IRNAargW</i>]	U11296	1.2e-21				
				Bacteriophage P4 integrase	X05947	1.8e-19				
				<i>E. coli</i> K12 integrase [<i>IRNAleuX</i>]	U14003	5.7e-19				
				<i>Salmonella typhimurium</i> LT2 Gifsy-1 prophage [<i>lepA</i>]	AF001386	1.1e-14				
				<i>Enetrobacter aerogenes</i> integrase	AF039582	4.2e-12				
				<i>Pseudomonas putida</i> integrase from the <i>clc</i> element [<i>IRNagly</i>]	PPAJ4950	1.2e-10				
				<i>Rhodobacter capsulatus</i> putative prophage integrase	U57682	3.2e-07				
				<i>Salmonella enteridis</i> integrase from IS3-like element	SEJ002209	1.4e-04				
				Bacteriophage adh site-specific recombinase IntG	Z97974	2.7e-02				
				Lambda Bacteriophage phi80 integrase	X04051	5.4e-02				
				<i>Streptococcus pyogenes</i> phage T12 integrase [<i>IRNAser</i>]	U40453	2.3e-01				
				<i>S. typhimurium</i> site-specific recombinase XerC	U92525	3e-01				
				regA	1793-2491	233	39.2/232aa	<i>Pseudomonas aeruginosa</i> genes for negative regulator of pyocin genes, PrtR	D12706	3e-33
								<i>E. carotovora</i> RgdA DNA binding protein	L32173	8.3e-32
Bacteriophage phi80 early region, <i>cl</i> gene product	X13065	1.4e-08								
Bacteriophage D3112 unidentified orf from <i>P. aeruginosa</i>	X52258	3.4e-08								
Bacteriophage 434 <i>cl</i> gene product	Y00118	6e-08								
<i>Staphylococcus aureus</i> Bacteriophage phi PVL proviral protein encoded by <i>orf19</i>	AB009866	3.9e-03								
<i>Haemophilus influenzae</i> unidentified putative protein	U32825	1.4e-01								
<i>Bacteroides thetaiotomicron</i> conjugative transposon <i>rteC</i> gene putative protein	L02419	3.0e-03								
<i>Mycoplasma bovis</i> PG-45 variable surface protein A (VspA)	L81118	9.4								
<i>Pseudomonas denitrificans</i> <i>orf7</i> putative protein downstream of <i>cob</i> genes	M62866	9.8e-55								
gepA	3526-3056	157	31.7/70aa	NSH ^b						
							22.0/127aa			
gepB	3624-4220	199	64.9/188aa	NSH ^b						
gepC	>4595-4798	>67	-	NSH ^b						

a. *intB* element from *D. nodosus* strain A198 (GenBank accession number X98546).

b. NSH indicates that there is no significant homology to sequences in databases.

c. Where the site of integration is known for an integrase gene it is indicated in square brackets following the description of the Int homologue.

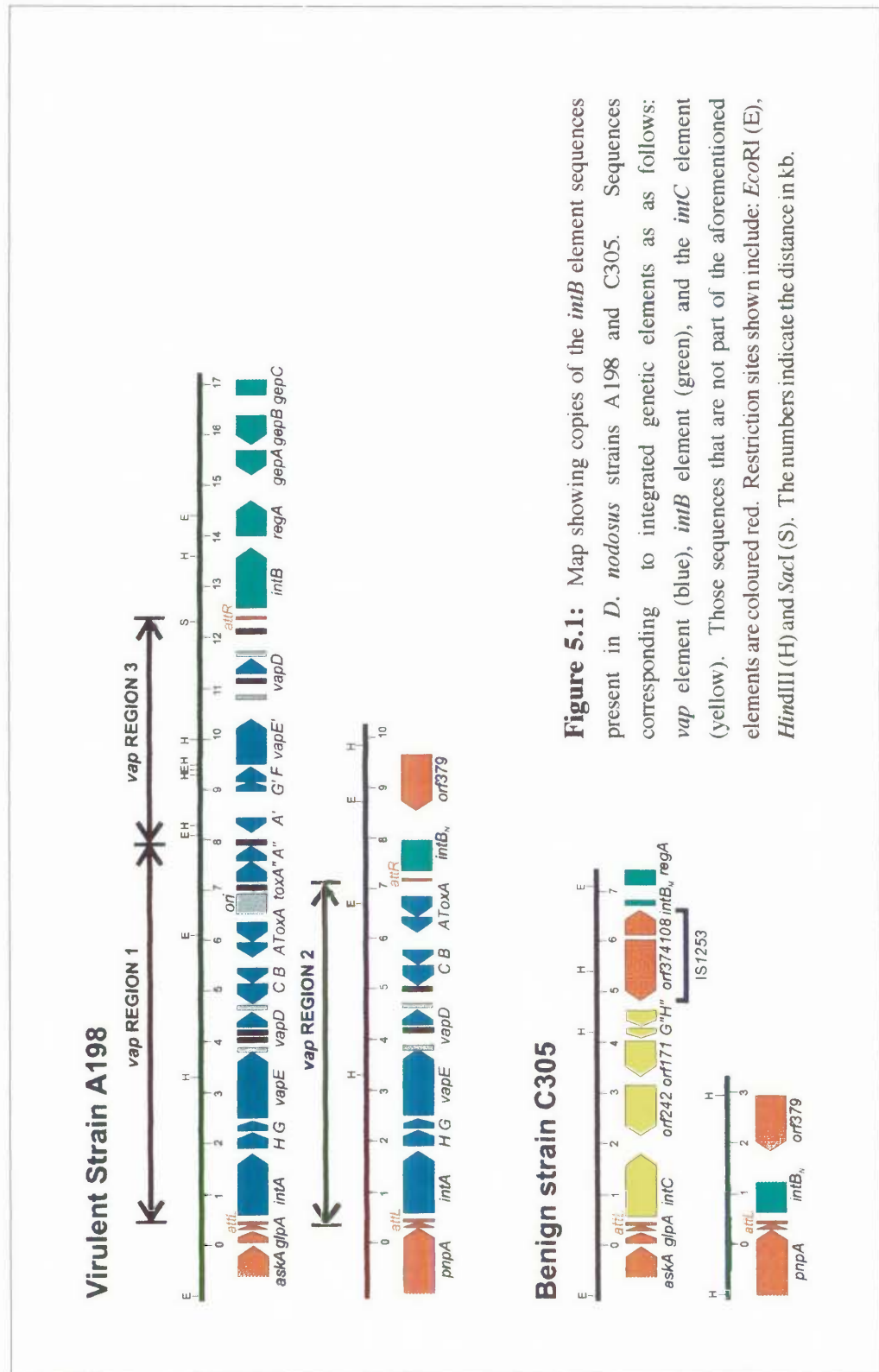


Figure 5.1: Map showing copies of the *intB* element sequences present in *D. nodosus* strains A198 and C305. Sequences corresponding to integrated genetic elements as follows: *vap* element (blue), *intB* element (green), and the *intC* element (yellow). Those sequences that are not part of the aforementioned elements are coloured red. Restriction sites shown include: *EcoRI* (E), *HindIII* (H) and *SacI* (S). The numbers indicate the distance in kb.

coding region, and consequently is herein called *intB_M* (Figure 5.1) (Bloomfield, 1997). *intB_M* is followed by *regA*.

In summary, in strains A198 and C305 three variants of the *intB* gene have so far been sequenced. In an effort to further characterise the *intB* element, chromosome walking to the right of *gepC* in strain A193 was undertaken, and a sequence of 4.3 kb determined. The prevalence, arrangement and integrity of the *intB* element in seventeen strains of *D. nodosus* was investigated in Southern blot experiments, in order to determine if there is a correlation between the presence and integrity of the *intB* element, and virulence.

5.2 Results

5.2.1 Screening of *D. nodosus* strain A198 lambda library

A library of genomic DNA from the virulent *D. nodosus* strain A198 in bacteriophage lambda was screened using DNA fragments (Figure 5.2) derived from the *intB* gene (probe 1) and *gepBC* (probe 2) to probe filters from duplicate plaque lifts. Four lambda clones (λ GW1.7, λ GW1.5, λ GW1.2 and λ GW1.4) hybridising to both probes and a fifth clone hybridising to probe 2 alone (λ GW1.6) were isolated (Figure 5.2). Restriction maps of overlapping lambda clones were constructed and aligned with previously-determined *intB* element sequences (Bloomfield *et al.*, 1997b) from strain A198 (Figure 5.2) which are adjacent to *vap* region 3 (Figure 5.1).

5.2.2 Sequencing of the region downstream of *gepC* in *D. nodosus* strain A198

DNA fragments from *EcoRI*, *HindIII* and *EcoRI/HindIII* restriction enzyme digests of λ GW1.4 were gel purified and subcloned. In an effort to confirm that the sequences

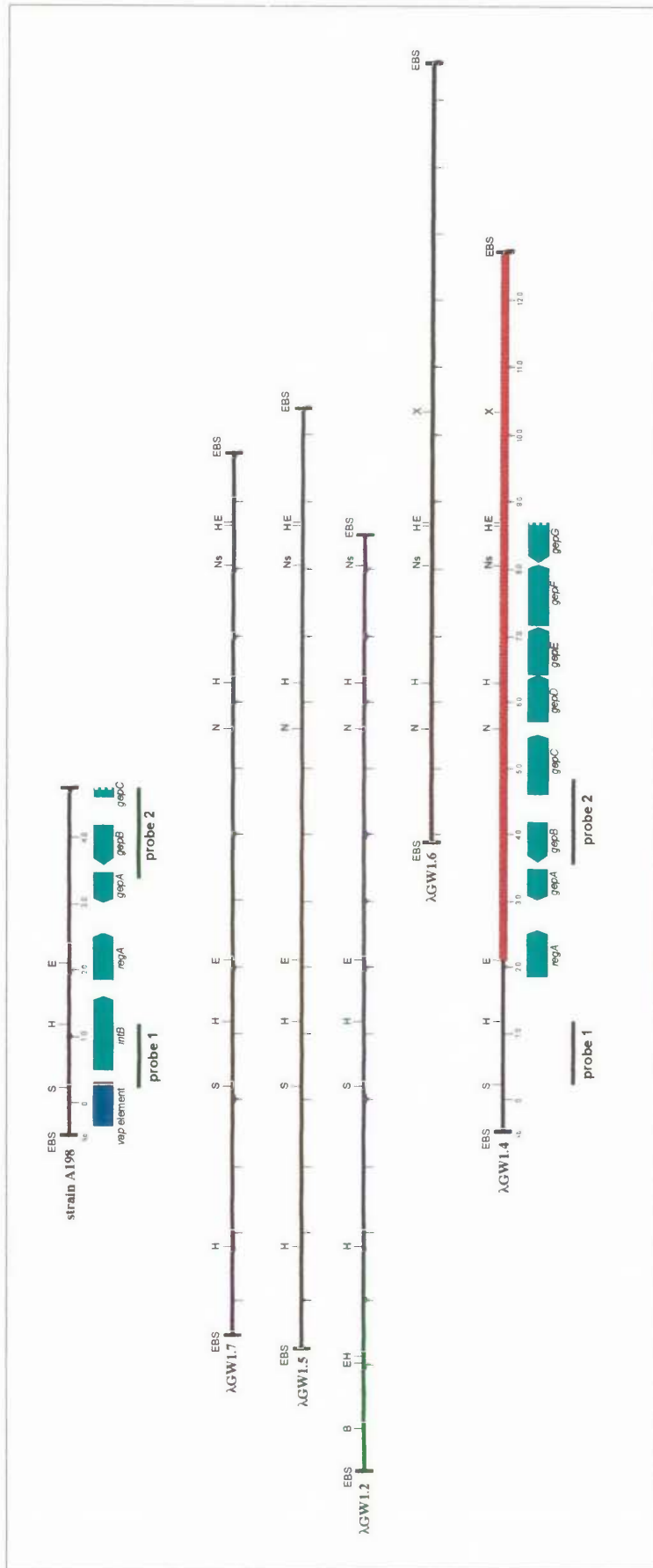


Figure 5.2: Restriction map of lambda clones isolated from a genomic library of *D. nodosus* strain A198 (Katz *et al.*, 1994), using probes specific for *intB* (probe 1) and *gepA-C* (probe 2). Restriction enzyme sites shown include: *Bam*HI (B), *Eco*RI (E), *Hind*III (H), *Nru*I (N), *Nsi*I (Ns) and *Sac*I (S). The numbers indicate the distance in kb. Sequences that were subcloned from λ GW1.4 are indicated by a red box.

isolated on λ GW1.4 reflected the arrangement of the same sequences in the A198 genome, a Southern blot analysis was performed. Genomic DNA from *D. nodosus* strain A198 was digested with *EcoRI*, *HindIII* and *EcoRI/HindIII* and probed with pGW39.1 (Figure 5.3). pGW39.1 hybridised to 7.0 kb *EcoRI* fragment, 5.3 kb and 2.5 kb *HindIII* fragments, and 4.4 kb and 2.5 kb *HindIII/EcoRI* fragments, which is consistent with the sizes predicted from the restriction enzyme map of λ GW1.4 (Figure 5.3). A series of subclones extending from the *EcoRI* site at position 2114 to the most distal multiple cloning site in the right arm of bacteriophage lambda were constructed (Figure 5.3). A sequence of 4261 bp was determined (from 4671 nt right of the *SacI* site to the *EcoRI* site at position 8932, Figures 5.3 and 5.4).

The *intB* element sequence (Figure 5.4) has a G + C content of 43%, which is similar to the 45% G + C content characteristic of the *D. nodosus* chromosome (Holdeman, Kelley & Moore, 1984). It was previously proposed that the *intB* element may have been acquired by horizontal transfer (Bloomfield *et al.*, 1997b), however, elements that have been acquired horizontally often have quite a different G + C content from the host chromosome. If the *intB* element was acquired horizontally, it may have been acquired long ago, and codon usage may have since evolved to be more like the host chromosome. Alternatively, it may have been acquired from an organism with a similar G + C composition.

The sequences upstream of *intB*, *regA* and *gepA-B* were analysed for open reading frames with start codons preceded by Shine-Dalgarno ribosome binding sites (Shine & Dalgarno, 1974), which led to the identification of five putative open reading frames, that have been designated genetic element protein genes *gepC*, *gepD*, *gepE*, *gepF* and *gepG* respectively in order to be consistent with earlier nomenclature (Bloomfield *et al.*, 1997b) (Figure 5.5).

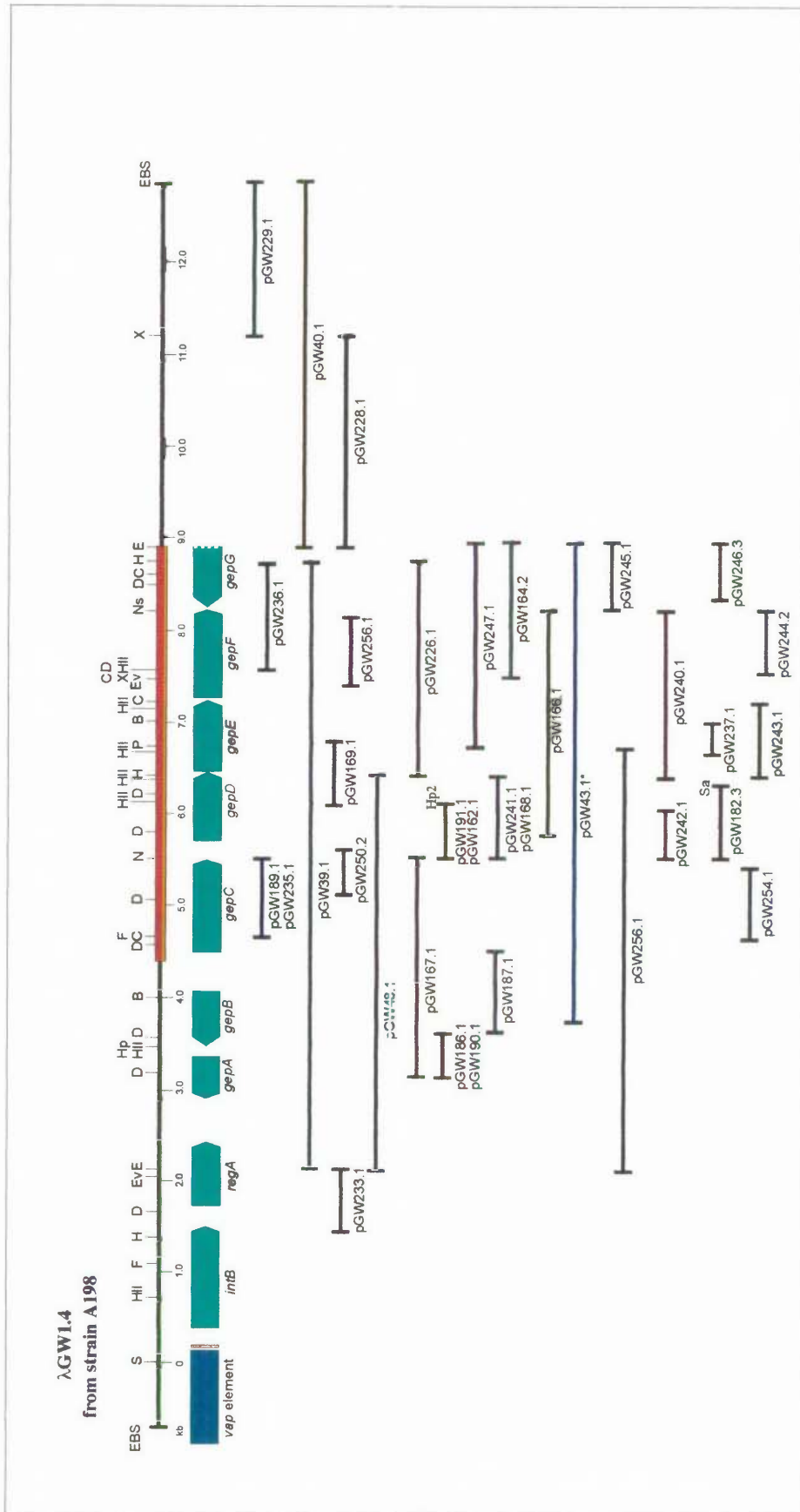


Figure 5.3: Restriction map of subclones from *D. nodosus* strain A198 lambda clone λGW1.4. Restriction sites are as follows: *Bbu*I (B), *Clal* (C), *Dra*I (D), *Eco*RI (E), *Fsp*I (F), *Eco*RV (Ev), *Hind*II (HII), *Hind*III (H), *Nru*I (N), *Nsi*I (Ns), *Sac*I (S) and *Xba*I (X). Restriction enzymes that recognise 4 bp recognition sequences are not shown on the map, but are instead shown at the ends of the subclones where applicable, as *Hpa*II (*Hpa*2) and *Sau*3AI (*Sa*). Note that only restriction sites used in cloning are shown upstream of the *Eco*RI site at position 8925. The red box indicates the double stranded DNA sequence determined in this work. An asterisk marks a subclone which was derived from λGW1.6 DNA (Figure 5.3). The numbers indicate the distance in kb.

```

-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
-35 -10
AGCAAACGATTGCATTAATAATTTATTACCCTTGTATTATTAAAAAAATCAAGCGCCAGCAAACGTTTCACCGCTGTTTT 4400
TCGTTTGCTAACGTAATTTTAATAATGCGGAACAAAATAATATTTTAGTTTCGGGTCGTTTGCAAAGTGGCGACAAAA

TTATCATACCAATCTGCAAAAATGGCAGCACATAAGTATTTTAAAACGATCATTGCTAAACTGTATTTCATTTGATGAA 4480
AATAGTAATGGTTAGACGTTTTTACCCTCGTGTATTTCAATAAATTTTGCTAGTAACGATTTGACATAAGTAAACTACTT

CTGCTGGTTTAAACAAAAGGAACAAGTGATGGCACACTTTTTTTATGGAACATTAATCCGCTGGCAACGCAGCGCGCAAA 4560
GACGACCAAAATGTTTTTCCCTTGTTCACTACCGTGGAAAAAAATACCTTGTAAATAGGGGACCGTTGCGTCCGCGGTTT

      gepC M P I D V F C F A K E C R V K A
AGGAACCGTAGAATTAGACGGCAGCAAATGCTTAAAGCCATTGATGTATTTTGTTTTGCTAAAGAATGCCGCGTAAAAAG 4640
TCCTTGGCATCTTAATCTGCCGTCGTTTACGAATACGGGTAACACATAAAACAAAACGATTTCTTACGGCGCATTTTC
  G A R L A F Q I D P K Y N T I F D A Q I V G F T Y S
CCGGCGCGCGCTTGGCATTTCAAATCGATCCCAAATACAACACCATCTTTGATCGGCAAATCGTTGGTTTTACCTACTCT 4720
GGCGCGCGCAACCGTAAAGTTTAGCTAGGGTTTATGTTGTGGTAGAAACTACCGGTTTAGCAACCAAAATGGATGAGA

F N S E Q E I D R F D D D L P S D F S L I I S A I F A
TTTAATTCGGAACAAGAAATTGATCGGTTTGATGAAGACTTGCCAGCGATTTTTCATTGATCATTCTGCCATTTTTC 4800
AAATTAAGGCTTGTCTTTAACTAGCCAAACTACTACTGAACGGGTCGCTAAAAAGTAACTAGTAAAGACCGTAAAAACG

A V F F A M I T F L Y P P S I L I L L G T A S V I M F
CGCCGATATTTTGGCATGATCACCTTTCTTTATCGCGCTATATATTGATTTTATAGGCACGGCTTCGGTCATCATGT 4880
CGCGATAAAAAACCGTACTAGTGAAAGAAATAGCGCGCAGATATAACTAAAAATAATCCGTGCCGAAGCCAGTAGTACA

  I Q L L Y E R W S W E R N Q S P L C S F N P P V F L
TTATTCAATTGCTTATGAACGCTGGTCGTGGGAGCGCAATCAATCGCCGTTATGTTCATTTAATCCGCCGTTTTTTTA 4960
AATAAGTTAACGAAACTTGGCAGCAGCACCCTCCGTTAGTTAGCGGCAATACAAGTAAATTAGCGGGCAAAAAAAT

I L S L L G A W S G A L L G Q Y M F N Y Q R R Q P R F
ATTTTATCATTGCTCGCGCGTGGTCAGGCGCTACTTGGGCAATATATGTTTAACTATCAACGGCGCCAACCGCGTTT 5040
TAAATAGTAACGAGCCGCGCACAGTCCGCGCAATGAACCCGTTATATACAAAATGATAGTTGCCCGGTTGGCGCAA

  K Y L L W L V S G I N F S V L F L L G I N F D M R P P
TAAATATTTATTGGCTCGTCTCCGGCTCAATTTCCGCTCCTATTTTATTAGGCATCAATTTTGACATGCGTCCGC 5120
ATTTATAAATAACACCGAGCAGAGCCGTAGTTAAAAAGGCAGGATAAAAAATAATCCGTAGTTAAAACTGTACGCAGCGC

  E P E P V A I V Q T H Q P S S L T E Q N T A P T D T
CGGAACCGAACCCGTCGCTATCGTGCAAAACCCATAACCGTCATCGCTGACCGAACAAAATACCGCGCCGACCGACACC 5200
GCCTTGGTCTTGGCGAGCAGATAGCACGTTTGGGTAATTGGCAGTAGCGACTGGCTTGTTTTATGGCGCGGCTGGCTGTGG

A P Q K P V T N T V S N Q T E P S S T D I S Q N S S P
GCACCGCAAAAACCCGTAACCAACTGTTCCAAICAAACCGAACCATCATCAACCGACATTTCGCAAAATTCATCGCC 5280
CGTGGCGTTTTTGGCATTGGTTATGACAAAGGTTAGTTTGGCTTGGTAGTAGTTGGCTGTAAAGCGTTTTAAGTAGCGG

  T A T A P D I S L P A V S F F E T P K V T P P E S C Y
AACTGCCACCGCGCGGACATTTCATTGCCTGCGGTTTCATTTTTGAAACACCAAAAGTTACGCCCGCGGAATCGTGTT 5360
TTGACGGTGGCGCGCCCTGTAAAGTAACGGACGCGCAAGTAAAAAAACTTTTGGGTTTCCAATGCGCGCGCCTTAGCACAA

  I V V Q D F R D M E A A K T F A A E Q A L N N P D A
ATATCGTCGTGCAAGATTTTCGTGATATGGAAGCACGAAAACTTTTGGCCCGCAACAAGCATTAACAATCTGATCGG 5440
TATAGCAGCAGCTTCTAAAGCACTATACCTTCGTCGCTTTTGAAAACCGCGGGCTTGTTCGTAATTTGTTAGGACTACGC
-35 -10
P I R I F F T Q K Y K F A V T N G T L E I S S A A A Q
CCAATCGCATCTTTTTACGCAAAAATAAATTCGCGCTACCAACGGCACATTAGAAATTTCATCTGCCCGCGCGCA 5520
GGTTAAGCGTAGAAAAGTGCTTTTTATATTTAAGCGGCAGTGGTTGGCGGTGAATCTTTAAAGTAGACGGCGCGCGT

  L D Q K I A A G E L P P E S Y C L L K A G V R E E V A
ATTAGATCAAAAATCGCCGCGGTGAATTGCCGCGGAAAGTTATTTGTTGCTGAAAGCCGGCGTTCGCGAAGAAAGTGG 5600
TAATCTAGTTTTTAGCGGCGGCCACTTAACGGCGCGCTTCAATAACAAACGACTTTCGGCCGCAAGCGCTTCTTCACC

R *
CGCGTAAAAAGATGGACAACCGTCCCGCGCGCTITTAATCTAAAAAAAACTTGTCCGGGTCGTTTTCAGCTGAGAAAT 5680
GCGCAATTTTCTACCTGTTTGCAGGCGCGCGCAAAATTAGATTTTTTGAACAGCCCCACGACAAAAGTCGACTCTTTA

-35
ACCCGTTGAACCTGATTAGTTAATACTGACGGAGGAAACAAGCCAACGCTTTTGTCTCCTCCGTTTTTGGAGGATTTTT 5760
TGGGCAACTTGGACTAAGTCAATTATGACTGCCTCCTTGTTCGGTTGCGAAACAGAGGAGGCAAAAACCTCCTAAAAA

-10      gepD M T S V T V S I K N
TTTGATGATTTTACTGTGCAATCCATCTGATTTTAAAGCGAGTAACGATGACTAGCGTTACCGTTTCCATCAAAAAT 5840
AACTACTAAAAATGACGAGCTTAGGTAGACTAAATTTTCGCTCATTGCTACTGATCGCAATGGCAAAGGTAGTTTTTA

```



```

-----+-----+-----+-----+-----+-----+-----+-----+-----+
L Y L S Y G A H V I F Q D F S H D F A A N A W H V I L      5920
CTTTATTTATCCTACGGCGCGCACGTCATTTTTCAGGACTTTTCCACGATTTGCGGCCAACGCGTGGCACGTCATTTT
GAAATAAATAGGATGCCGCGCGTGCAGTAAAAAGTCTGAAAAGGGTGCTAAAACGGCGGTTGCGCACCGTGCAGTAAAA

  G R S G C G K S S L L F A I A G L L A A N G Q Q R G S      6000
AGGACGCTCGGGCTGCGGTAAAAGCAGTTTATTACACGCCATTGCGCGGATTACTTGCTGCCAACGGAACAGCGCGGAA
TCCTGCGAGCCCGACGCCATTTTCGTCAAATAATCGCGGTAACGGCTAATGAACGACGGTTGCCTGTTGTCGCGCCTT

  I D D G A G N S L S G R L R W M A Q E N D L L P W L      6080
GCATTGATGACGGCGCTGGCAACTCATTATCCGGTGCATTGCGTTGGATGGCGCAGGAAAATGATTTATTGCCGTGGTTC
CGTAACTACTGCCGCGACCGTTGAGTAATAGGCCGCTAACGCAACCTACCGGTCCTTTTACTAAATAACGGCACCAAC

  N I E D N V L L G A H V R G E A K N A A A V E R L L T      6160
AATATTGAAGACAACGTTCTTTTGGGCGCGCACGTCGCGCGTGAAGCGAAAATGCGGCGGCGAGTGGAGCGTTTGTTCAC
TTATAACTTCTGTTGCAAGAAAACCCGCGCGTGCACGCGCCACTTCGCTTTTACGCGCGCGTACCTCGCAAAACAATG

  A C G L K V N N K K R F Q Q L S G G E R Q R L A L A R      6240
GGCTTGGCGTTTAAAAGTCAACAACAAAAGCGTCGCGAGCAATTATCCGGCGGCGAACGGCAACGGCTGGCTTTAGCGC
CCGAACGCCAAATTTTCAGTTGTTGTTTTTCGCAGCGCTGTTAATAGGCCGCCGCTTGCCTTGCAGCCGAAATCGCG

  T L I D D A P L I L M D E P F S A L D A I T R Y Q L      6320
GCATTTAATTGACGATGCCCGCTCATTGTTGATGACGCAACCTTTTCCGCTTTAGACCGGATACGCGTTATCAATTG
CGTGAATTAACGCTACGGGCGAGTAAAACCTACGCTTGAAAAGGCGAAAATCTGCGCTAGTGGCAATAGTTAAC

  Q N L A T E L L V D R T V I M I T H D P A E A L R L A      6400
CAAAATCTCGCACCGAATTATTGTTGACCGCACCGTGAATTATGATCACGCGATCCGGCAGAACTTTGCGCTTAGC
GTTTTAGAGCGGTGGCTTAATAACCAACCGCGTGCACATAACTAGTGCCTGCTAGGCCCTTTCGAAAACGCGAATCG

  N Y L Y V L E N G A L I E L P L P A A A P P R A F T S      6480
AAATATTGTACGTTTTAGAAAACGGCGCGTACCGAATTACCGCTTCCGCGCGCGCACCGCGCGGCTTTACCA
TTTTAATAACATGCAAAATCTTTTCCGCGCGATGGCTTAATGGCGAAGGGCGGGCGTGGCGCGCGCGCAATGGT

                                     gepE M I R F A A Q F L L T A
  E G F A E R Q Q Q L L E H L R *
GCGAAGGTTTTGCCGAACGGCAACAGCAACTTTTGGAGCATTTGCGATGATTCGTTTTGCCGCACAATTTTTATTAACG 6560
CGCTTCCAAAACGGCTTGCCTTGTGTTGAAAACCTCGTAAACGCTACTAAGCAAAACGGCGTGTAAAAATAATTGAC

  F G L C F L W Q S V I W L T Q T P P Y I L P S P L A      6640
CATTCGGACTGTGTTTTTATGGCAAAGCGTCATTTGGTTAACGCAGACGCGCCCTTATATTTGCCGTGCGCGCTGGCA
GTAAGCCTGACACAAAAATACCGTTTTCGAGTAAACCAATTGCGTCTGCGCGGAATATAAAAACGGCAGCGCGACCGT

  V L Q C L Y T D F D V L W R H S R V T V L E I A L S L      6720
GTTTTGCAGTGTATACACAGATTTGATGTTTATGGCGCCACAGCGCGTTACCGTATTAGAAATGCTTTGAGTTT
CAAAACGTCACAAATATGTCTAAAACCTACAAAATACCGCGGTGTCGGCGCAATGCATAATCTTTAACGAAACTCAA

  G I A S V F G T A T A I I L M I N A R L R R W L Q P L      6800
GGGATTCGCGCGTTTTTGAACGGCAACGGCGCTTATTTAATGATTAATGCGCGTTTGGACGATGGCTGCAGCCAT
CCCCTAACGGTTCGCAAAACCTTGCCTTGCCTGCAATAAAATTAATAATACGCGCAAACTCTGCTACCGACGTCGGTA

  I L V S Q T M P V Y A L A P I L M L W F G Y G L T P      6880
TAATTTTATGATCGCAACGATGCCCGTTTACGCCCTCGCACCAATTTTAAATGCTTTGGTTTGGTTACGGTTTAAACGCG
ATTAATAATCATAGCGTTTGTACGGGCAATGCGCGAGCGTGGTTAAAATACGAAACCAACCAATGCCAAATGCGGCG

  K I I V T V L I A Y F P I T T A T F D G L Q Q T P P A      6960
AAAATCATCGTTACCGTGTAAATCGCTTACTTTCCATACCACCGCCACTTTTGGCGTTTACAACAAACGCGCGCGC
TTTTAGTAGCAATGGCACAATTAGCAATGAAAGGTAGTGGTGGCGGTGAAAACGCAAAATGTTGTTTGGCGCGCGC

  Y L R L A Q T L G A N E R Q I L W R I R M P A A L P H      7040
TTATTTGCGTTTACGCAACGTTAGGCGCGAATCGCGGCAAAATTTATGGCGCATTCGCATGCCGCGAGCGTTGCCGC
AATAAACGCAAAATCGCTTGAATCCGCGCTTACGCGCGTTTAAAATACCGCGTAAGCGTACGGCGCTCGCAACGGCG

  L A S G L R V G A A M A P I G A I I G E Y V G G S D      7120
ATTTGGCATCGGGATTGCGCGTGGGCGGGCGATGCGCCGATTGGTGCCATTATTTGGTGAATATGTTGGCGGAAGCGAT
TAAACCGTAGCCCTAACGCGCACCCGCGCGCTACGCGGCTAACACGGTAATAACCACTTATACACCCGCTTCGCTA

  G L G Y L M Q Y G I N R S Q V A L T F A A L F V M T L      7200
GGGTGGGTTATTTAATGCAATACGGCATTAATCCAGCCAAGTTGCGTTAACTTTTGCCTTTGTTTGTGATACGTT
CCCAACCAATAAATACGTTATGCGGTAATAGCTCGGTTCAACGCAATGAAAACGGCGAAACCAACACTACTGCAA

  L T L A I Y Y G I D A I F E K M V L L G N N G E *gepF M      7280
ATTAACGTTAGCGATTTATTACGGCATCGATGACTTTTTGAAAATGGTATTATTGGTAAACAATGGGGAGTAATACA
TAATTGCAATCGCTAAAATAATGCCGTAGCTACGTCAAAACCTTTTTTACCATAATAACCAATGTTTACCCCTATTATGT

```

```

-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
K K I S T F L F G L M L A T T A L A K E P L H L M L
TGAAAAAATCAGTACGTTTTTATTCGGATTAATCTTGGCAACAACGGCTTGGCAAAAGAGCCGCTGCATTTAATGTG 7360
ACTTTTTTTAGTCATGCAAAAATAAGCCTAATTAC AACCGTTGTTGCCGAAACCGTTTTCTCGGCGACGTAATTAACAAC

D W F I N P N H A P I I I A Q Q N G Y F D K H G L E V
GATTGGTTTTATCAATCCCAATCATGCGCCGATCAI CATCGCGCAACAAAACGGTTATTTTGATAAACACGGTTTAGAAGT 7440
CTAACCAAAATAGTTAGGGTTAGTACGCGGCTAGTAGTAGCGGTTGTTTTGCCAATAAAACTATTTGTGCCAAATCTTCA

T I T E P S D P A L P F K L V A A E K I D L A I N Y Q
CACCATCACCGAACCGTCCGATCCGGCGCTGCCGCGAAATGGTTGCCGCGGAAAAAATCGATTTAGCCATCAATTATC 7520
GTGTTAGTGGCTTGGCAGGCTAGGCCGCGACGGCGCTTTAACCAACGGCGCTTTTTTAGCTAAATCGGTAGTTAATAG

Q Q L H L Q I D E G L P I S R V S T L I A T P L N C
AACAACAATTCGATTTACAATTTGACGAAGGATTC CGGATATCGCGCGTATCAACGTTAATCGCTACGCCCTTAAATTCG 7600
TTGTTGTTAACGTAATGTTTAACTGCTTCTCAACGGCTATAGCGCGCATAGTTGCAATTAGCGATCGGAAATTTAACG

V I V D A Q S G I K Q V S D L K G K K I G Y S V A G V
GTGATTGTTGACGCGCAAAAGTGGCATTAAACAAGT TCTGTATTTAAAGGTTAAAAAATGGTTATTTCCGTTGCCGGCGT 7680
CACTAACAACTGCGCGTTTTACCGTAATTTGTTCAAGACTAAATTTCCATTTTTTTAACCAATAAGGCAACGGCCGCA

D E A V L Q S F L A S G G L T L N D V K L V N V N F S
TGATGAAGCCGTGTGCAATCTTTTTTAGCCAGCCGTGGTTAACTTTAAACGATGTGAAACTCGTCAACGTCATTTTTT 7760
ACTACTTCGGCACAACTTAGAAAAAATCGTCCG CACCAATTGAATTTGCTACACTTTGAGCAGTTGCAGTTAAAAA

L S P A V M S G Q V D A V I G A A R T V E L H E M K
CTTTATCGCCGGGTGATGAGTGGACAAGTAGAC GCGGTTATCGGCGCCGCGCACGGTGAATTACACGAAATGAAA 7840
GAAATAGCGCCGCCACTACTACTCCTGTTTCACTGCGCCAATAGCCGCGCGCGTCCACCTTAATGTGCTTTACTTTT

A H N H E G R A F F L E E H G I P P F D E L I F V A H
GCGCAATCATGAAGGGCGCGCTTCTTTTTAGAGAACAACGCGCATCCGCCCTTGTATGAATGATTTTTGTGCGCGCA 7920
CGCGTGTAGTACTTCCCGCGCGGAAGAAAAATCTCTTTGTGCCGTAAGGCGGAAACTACTTAACAAAAACAGCGCGT

N K H R H D E K I V K F N E A L T E A V H F I V N H P
TAACAAACACCGCCACGATGAAAAAATCGTTAAAT TTAATGAGGCGCTCACCGAAGCCGTGCATTTTTATTGTCAATACC 8000
ATTTGTTGTGGCGGTGCTACTTTTTTAGCAATTTAAATTACTCCGCGAGTGGCTTCGGCACGTAATAACAGTTAGTGG

E E A W Q K Y I A Y K K G L D D A V N Q Q A W K D S
CTGAAGAAGCATGGCAAAAATATATGCTTTATAAAAAGGTTTAGATGATGCGGTTAATCAACAAGCATGGAAAGACAGT 8080
GACTTCTTCGTACCGTTTTTATATAACGAATATTT TTTCCAAATCTACTACGCCAATTAGTTGTTTCGTACCTTTCTGTCA

L T R F A L R P A A L D D R R Y Q N Y A Q Y L H Q I G
TTAACGCGTTTTGCTTTGCGTCCGCGCGGCTTGA TGATCGGCGTTATCAAAATTACGCGCAATATTTGCATCAAAATGG 8160
AATTGCGGAAACGAAACGCGAGCCGCGCAACTACTAGCCGCAATAGTTTAAATGCGCGTTATAAACGTTAGTTTAAAC

L I K K I V P V S E Y A V Q P *
TTTAATTAATAAATCGTGCCGGTTTCGGAGTATG CCGTGCAACCGTAATCGCGTTAATCCTAAAAAGCCCTGTAACAG 8240
AAATTAATTTTTTTAGCACGGCCAAAGCCTCATAC GGCAGTTGGCATTAGCGCCAATTAGGATTTTTTCGGGACATGTC

TGCTTTTTTATGCATTTTATTGGCACGGTTATAA CAAAAATGCCCGTGAGCTTTTTTATATAAAGAAACGGCGTAAC 8320
ACCGAAAAAATACGTAAAAATAACCGTGCCCAAT ATGCTTTACGGGCACTCGAAAAATATATTTCTTTGCCCGCATGTAT
* I F F R A Y S

TCGTCATGCCGCTGACATAATCGCACACCGCCAA TAATTTTTGATACGGGGAATCCGCATAACGACGATAACCGTTGG 8400
AGGCAGTACGGCGACTGTATTAGCGTGTGGCGGT TATTAATAACTATGCCCTTAGGCGTATTGCTGCTATTGGCGAACC
D T M G S V Y D C V A L L K Q Y P S D A Y R R Y R T P

CAATAATAATAAATCATTTTCATCGTGCCGGCGGC IGGTTTTATTGTTTAAAAACCGTATTGGCAACCGCGTGCAAAAAA 8480
GTTATTATTATTTTAGTAAAGTAGCACGGCCGCGA CCAAATAACAAATTTTGGCATAACCGTTGGCGCCACGTTTTTT
L L L L I M E D H R R S T K N N L V T N A V A T C F

TATCCAATAATCCGCGCAAAATGCGGTTCCCGCC ACTTGAATTTCCACAACCGGTTGCAAGTTATAAATAAATCGCTG 8560
ATAGGTTATTAGGCGCGTTTTACGCCAAGGGCGG TGAACCTAAAGGTGTTGGCCAACGTTCAATATTTATATTAGCGAC
I D L L G R L I R N G A V Q I E V V P Q L N Y I Y D S

GCAAATGGCGCAACGCTTTTAAATGACGATAAGA CCGAATATGCGGCAATAAAGCAGCGCGGTTCCATTTAAAAATTCG 8640
CGTTAACCGCGTTGCGAAAAATTAAGTCTATTCT CCGTTATACCGGTTATTTCTGTCGCGGCAAAAGGTAATTTTAAACG
A F Q R L A K L Q R Y S P I H P L L A A G N G N L I A

CGATTCTGATTGCCAAATTTGCGCGCACGTTTCA TCGATTAATTTGCCAATCGCTTTGGCGCGCAATAACCCAATTTAT 8720
GCTAAGACTAACGTTTTAACGCGCGTGCAAGTA 3CTAATTAACAGGTTAGCGAAACCGCGGTTTTATTGGGTTAAATA
S E S Q W F Q A C T E D I L Q G I A K A R L Y G L K

```

```

-----+-----+-----+-----+-----+-----+-----+-----+
CGCTTTTACTGTGCATTGCTTTGATGCGCGTTTCATCAACATTTTCCCGCACCAATTCAAATAAACGCGCATAAGCTTCT 8800
GCGAAAATGACACGTAACGAAACTACGCGCAAAGTAGTTGTAAAAGGGCGTGGTTAAGTTTATTTGCGCGTATTCGAAGA
D S K S H M A K I R T E D V N E R V L E F L R A Y A E

TCAAATGAATTTGTTTCACCTGATAAGCGTCTTCAATATCAACGATTAATAACAATATCATCAGCCGCTTCCATCAA 8880
AGTTTACTTAAACAAAGTGGACTATTCGCAGAAGTTATAGTTGCTAATTTATTGTTTATAGTAGTCGGCGAAGGTAGTT
E F H I Q K V Q Y A D E I D V I L Y C I D D A A E M L

ATACGATAACGGGTGCCGAAGCCAAACTTGGCGTTTAATTTGCGGAATTC 8930
TATGCTATTGCCACGGCTTCGGTTTGAACGCCAATTAACGCCTTAAG
Y S L P H R L W V Q P K I Q P I gepG

```

Figure 5.4: Nucleotide sequence of sequences upstream from *intB*, *regA* and *gepA-B* from *D. nodosus* strain A198. The amino acid sequences of five putative proteins encoded by *gepC*, *gepD*, *gepE*, *gepF* and *gepG* are aligned with the nucleotide sequence. Shine-Dalgarno sequences are indicated in green type, and putative -35 and -10 promoter sequences (red type).

5.2.3 Identification of a putative binding protein-dependent ABC transporter

The first 67 aa of the putative *gepC*-encoded protein had been determined previously (Bloomfield *et al.*, 1997). In this work the remaining 270 aa of the putative GepC protein were determined. GepC has no similarity to sequences in the GenBank databases.

gepD encodes a potential protein of 239 aa which has a very high degree of amino acid identity to numerous ABC (ATP-binding Cassette) transporter, ATP-binding proteins (Table 5.2). *gepD* is in an operon-like arrangement with *gepE*, which overlaps the *gepD* coding region by one nucleotide. *gepE* encodes a protein which has very high similarity to ABC transporter membrane-associated proteins (Table 5.2). *gepF* encodes a protein that shares significant amino acid identity to Nmt1-like thiamine biosynthesis proteins (Table 5.2), which are, at least in *Haemophilus influenzae*, located immediately adjacent to genes encoding an ABC-transporter ATP-binding subunit and an ABC-transporter membrane-associated protein respectively (Table 5.2). *gepG* encodes a protein that shares similarity to hypothetical proteins of unknown functions.

Table 5.2: Sequence analysis of *gepC* to *gepG* from the *intB* element^a in *D. nodosus* strain A198

Gene	Co-ordinates 5'-3- (nt)	Size (aa)	% identity to nt or putative protein	Homologues/Description	Accession Number	P(n)
<i>gepC</i>	4595-5606	337	-	NSH	-	-
<i>gepD</i> ^b	5811-6528	239	40.8/240	<i>Haemophilus influenzae</i> HI0354 ABC transporter, ATP-binding protein	U32720	9.2e-24
			38.1/218	<i>Escherichia coli</i> K12 MG1655 Ort255 ABC transporter, ATP-binding protein [taurine]	D85613	9.0e-16
			30.9/239	<i>Archaeoglobus fulgidus</i> unidentified orf	U73857	4.3e-15
			34.2/242	<i>Pseudomonas putida</i> SsuB, ABC-type transporter, ATP-binding protein [sulfate]	AF075709	1.2e-14
			35.9/223	<i>Pseudomonas aeruginosa</i> AtsC, ABC-type transporter, ATP-binding protein [sulfate]	Z48540	1.6e-14
			36.1/202	<i>A. fulgidus</i> sulfate ABC transporter, ATP-binding protein AF0092 [sulfate]	AE001100	3.4e-14
			33.1/241	<i>Bacillus subtilis</i> YgaL hypothetical protein near nitrate transporter [nitrate]	Z93102	5.4e-14
			39.4/188	<i>B. subtilis</i> HisP importer [Histidine]	AF008220	6.2e-14
			34.5/246	<i>Mycobacterium smegatis</i> L5 and D29 bacteriophage resistance protein	U50335	8.2e-14
			30.1/229	<i>Serpulina hydroxybenzoate</i> ShiC putative ABC transporter [iron]	SHU75349	9.2e-14
			36.9/198	<i>Phormidium laminosum</i> (cyanobacterium) nrtC-PhI ATP-binding protein [nitrate]	Z19598	2.1e-13
			35.3/221	<i>Rhodobacter capsulatus</i> SB1003 ABC transporter, ATP-binding protein	AF010496	2.3e-13
			33.6/220	<i>H. influenzae</i> HitC iron utilisation protein [iron]	S72674	3.7e-13
			33.1/241	<i>B. subtilis</i> YglA, homologous to nitrate ABC transporter, ATP-binding protein [nitrate]	Z99108	4.7e-13
			34.6/205	<i>Synechococcus</i> sp. <i>ntrD</i> encoded nitrate transporter [nitrate]	X74597	8.7e-13
			36.4/217	<i>Methanococcus jannaschii</i> M0409 hypothetical protein	U67490	1.1e-12
			35.6/199	<i>Helicobacter pylori</i> HP0818 osmoprotection protein	AE000593	2.1e-12
			32.4/219	<i>Aspergillus nidulans</i> sulfate permease encoded by <i>cysA</i> [sulfate]	J04512	2.7e-12
			32.7/219	<i>Bacillus firmus</i> NatC ABC transporter, ATP-binding protein	AF084104	2.8e-12
			32.2/214	<i>Streptomyces coelicolor</i> GabT aminotransferase	AL031225	2.9e-12
			35.3/201	<i>Klebsiella pneumoniae</i> nitrate NasD nitrate transporter [nitrate]	KPNNASFEC	9.9e-12
			34.5/194	<i>H. influenzae</i> Rd thiamine ABC transporter, ATP-binding protein [thiamine]	U32782	1.0e-11
			31.2/231	<i>Salmonella typhimurium</i> ATPase of 2-aminoethylphosphonate transporter	U69493	1.5e-11
			30.2/248	<i>Agrobacterium tumefaciens</i> pTi15955 MotB ATP-mannopine transport protein	U60011	1.7e-11
			32.3/211	<i>Streptococcus mutans</i> MsmK mannose transport protein [mannose]	M77351	2.3e-11

Table 5.2 continued: Sequence analysis of *gepC* to *gepG* of the *intB* element^a in *D. nodosus* strain A198

<i>gepE</i>	6527-7274	249	<i>H. influenzae</i> HI0355 ABC transporter permease protein [unknown]	U32720	2.8e-49
			<i>Chlamydia trachomatis</i> ABC transporter permease CT854 [unknown]	AE001358	4.4e-23
			<i>E. coli</i> K12 MG1655 TauC [taurine]	AE000143	1.2e-21
			<i>B. subtilis</i> YgaM, homologous to nitrate permease protein [nitrate]	Z93102	1.9e-17
			<i>P. aeruginosa</i> <i>aisB</i> encoded ABC-type transporter, membrane subunit [sulfate]	Z48540	5.3e-16
			<i>M. jannaschii</i> MJ0409 hypothetical protein	U67490	1.8e-15
			<i>A. fulgidus</i> sulfate ABC transporter permease protein AF0093 [sulfate]	AE001100	7.7e-15
			<i>E. coli</i> K12 MG1655 OmpF outer membrane protein F precursor	AE000195	7.0e-14
			<i>E. coli</i> K12 ProW from osmoregulation operon	M24656	2.4e-13
			<i>Helicobacter pylori</i> HP0818 osmoprotection protein	AE000593	3.7e-13
			<i>P. aeruginosa</i> SsuC ABC transporter membrane subunit [sulfate]	AF075709	1.2e-12
			<i>E. coli</i> ProW glycine, betaine, proline transport protein	D90891	1.3e-12
			<i>B. subtilis</i> ProX osmoprotection protein	U38418	1.5e-10
			<i>Synechococcus</i> CmpB, homologous to NrtB [nitrate]	D26358	3.6e-10
			<i>A. fulgidus</i> AF0989 hypothetical protein	AE001036	8.2e-10
			<i>Synechocystis</i> sp. PCC6803 sIII715 hypothetical protein	D90916	1.1e-09
			<i>B. subtilis</i> choline transporter OpuBD	AF008930	1.5e-09
			<i>Synechococcus</i> PCC7942 periplasmic substrate CynB [cyanates]	AF001333	1.2e-08
			<i>B. subtilis</i> choline transporter OpuAB transmembrane protein	U17292	5.8e-08
			<i>B. subtilis</i> YvdK hypothetical protein	Z99121	6.8e-08
			<i>M. smegmatis</i> L5 and D29 bacteriophage resistance protein	U50335	8.6e-08
			<i>K. pneumoniae</i> NasE nitrate transporter [nitrate]	L27431	1.9e-07
			<i>E. coli</i> <i>dld</i> encoded lactate dehydrogenase	AE000302	6.4e-07
<i>Mycobacterium tuberculosis</i> H37RV DppB peptide transporter	AL022121	7.0e-05			
<i>Methylbacterium extorquens</i> <i>abcB</i> putative ABC transporter subunit B	U72662	6.0e-04			
<i>gepF</i>	7355-8207	309	<i>H. influenzae</i> HI0357 putative thiamine biosynthesis protein encoded by <i>nmtI</i>	U32720	3.5e-82
			<i>Ochrobactum antropi</i> hypothetical protein	AQ242226	1.4e-46
			<i>Uromyces fabae</i> Nmt1 thiamine biosynthesis protein	UFU8179	8.7e-17
			<i>Saccharomyces pombe</i> Nmt1 thiamine biosynthesis protein (chromosome III)	J05493	1.3e-16
			<i>Aspergillus parasiticus</i> Nmt1	U15196	5.1e-16
			<i>Saccharomyces cerevisiae</i> THI5 thiamine biosynthesis protein (chromosome XIV)	AL031579	4.3e-15
			<i>S. cerevisiae</i> Nmt1 thiamine biosynthesis protein (chromosome IV)	Z74292	4.8e-15

Table 5.2 continued: Sequence analysis of *gpc* to *gpg* of the *intB* element^a in *D. nodosus* strain A198

<i>gpf</i> (cont.)	7355-8207	309	29.1/175	<i>S. cerevisiae</i> OrfY (chromosome X)	Z49656	1.4e-14
			22.4/232	<i>C. trachomatis</i> fumC encoded fumarate hydratase	AE001358	8.2e-11
			25.0/272	<i>A. fulgidus</i> AF0091 encoded protein gene, adjacent to sulfate transporter [sulfate]	AE001100	1.1e-10
			23.6/245	<i>B. subtilis</i> orfK encoded hypothetical protein	L16808	1.4e-05
			23.6/245	<i>B. subtilis</i> YzeA, homologous to nitrate transport protein precursor	Z93102	5.6e-05
			36.4/77	<i>P. aeruginosa</i> atxA encoded arylsulfatase	Z48540	6.6e-04
			28.9/149	<i>Yersinia pestis</i> plasmid pMT1 unidentified orf near ABC transporter	AF053947	1.0e-03
			20.3/256	<i>M. bryantii</i> copper response extracellular protein	U40213	3.3e-02
			27.5/149	<i>Synechococcus</i> PCC7942 periplasmic substrate binding protein CynA, NrtA-like	AF001333	9.5e-01
			33.0/203	<i>Synechocystis</i> sp. PCC6803 s110189 hypothetical protein	D64002	2.1e-17
			25.5/161	<i>H. influenzae</i> HI1299 hypothetical protein	U32809	1.0
<i>gpg</i>	>8930-8297	>210	26.8/119	<i>P. aeruginosa</i> unidentified orf upstream of <i>cysB</i> encoding a transcription factor	U95379	10

a. NSH indicates that there is no significant homology to sequences in databases;

b. Note that only the top twenty five homologues of these proteins have been shown.

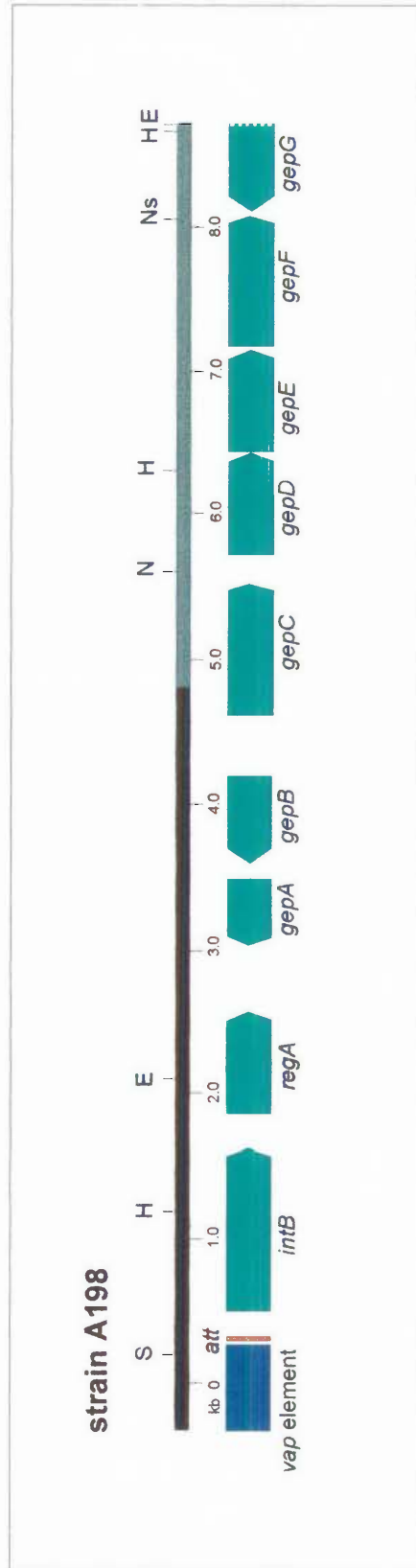


Figure 5.5: Restriction map of the *intB* element-associated sequences (green) from *D. nodosus* strain A198. The numbers shown indicate the distance in kb. Restriction sites shown include *EcoRI* (E), *HindIII* (H), *NruI* (N), *NsiI* (Ns) and *SacI* (S). The sequences determined prior to this work are indicated by a grey line, whilst those sequences determined in this work are indicated by a blue box. A putative attachment site, *att*, is also shown (red rectangle).

It is interesting that the proteins to which the GepD and GepE putative transport proteins are most similar are involved in the import of substrates such as sulfate, nitrate, histidine, iron, thiamine, mannose and glycine/betaine (Table 5.2), and all belong to the binding protein-dependent (BPD) ABC transporter subfamily (Boos & Lucht, 1996).

These bacterial binding protein-dependent (BPD) ABC transporters belong to a superfamily of proteins called ABC transporters, which share a region of high homology which extends over a region of 200 aa (Boos & Lucht, 1996; Higgins, 1992) in one of the polypeptides. This 200 aa region is called the ATP-binding cassette, within which there is a highly conserved ATP-binding motif that contains two conserved sites (A and B) that form an ATP-binding pocket (Rossmann *et al.*, 1975).

This ATP-binding site typically occurs at the end of an α -helix, and the amino acid consensus sequence forms a turn which brings the lysine (K) residue in close proximity with one of the phosphate groups in the Mg^{2+} -ATP (Fath & Kolter, 1993). The negatively charged aspartate residue (D) interacts with the positively charged Mg^{2+} ion in the A-site (Fath & Kolter, 1993). The consensus sequence for the A-site has been defined as [AG]-x4-GK-[ST] (Fath & Kolter, 1993), whilst the B-site is defined as VLLLDEP (Boos & Lucht, 1996). These binding sites have been subsequently identified in bacterial proteins that are dedicated to the export and import of substrates and other proteins that are dependent on cellular energy (Fath, 1993; Higgins *et al.*, 1986).

The ATP-binding cassette also contains another conserved sequence of 9 bp called the linker peptide that is thought to play a critical role in signal transduction between the ABC subunit and the integral membrane protein domains (Boos & Lucht, 1996).

Using the WAG motif program and the Prosite patterns database, an ATP-binding cassette was identified within GepD, and was aligned with the ABC transporter consensus sequence as described in the literature (Boos & Lucht, 1996; Fath & Kolter, 1993)

(Figure 5.6). The ABC transporter, ATP-binding protein is not involved in substrate specificity (Boos & Lucht, 1996) but hydrolyses ATP and thereby provides the energy for the translocation of a wide variety of substrates across membranes (Higgins, 1992). The linker peptide motif is also present in GepD (Figure 5.6).

The membrane-associated protein (permease) component of BPD ABC transporters has a strongly hydrophobic character. Computer aided analysis (Claros, 1996; Claros & von Heijne, 1994) of hydrophobicity of the GepE protein, and predicted topology suggests that the GepE protein is an integral membrane protein that contains four putative transmembrane domains (Figure 5.7). Furthermore, these putative transmembrane domains are separated by hydrophilic peptide loops of varying length, two exposed to the periplasmic face and one to the cytoplasmic face. Both the amino and carboxyl termini are located on the cytoplasmic face of the membrane according to the “positive-inside rule” (von Heijne, 1986).

Most of the membrane-associated proteins of BPD transporters that have been described previously are predicted to have six transmembrane domains that form three periplasmic loops and two cytoplasmic loops, with both C- and N- termini on the cytoplasmic face (Higgins, 1992). However, few BPD ABC importers have been characterised (Boos & Lucht, 1996), and of those that have been characterised, some variation in the number of transmembrane domains has been noted, as determined by gene fusion and proteolysis experiments. For example, the ProW (514 aa) protein of *E. coli* contains seven transmembrane domains (Hennessey & Broomesmith, 1993), whilst HisQ (228 aa) and HisM (235 aa) of *S. typhimurium* are short for proteins in this class (~300 aa) and have only five transmembrane domains (Traxler, Boyd & Beckwith, 1993). The *D. nodosus* GepE (249 aa) protein is similarly short, and thus may not conform to the six transmembrane-spanner consensus structure characteristic of ABC transporters in general



Figure 5.6: ClustalW alignment of ABC transporter, ATP-binding proteins with similarity to the *gepD* encoded ATP-binding cassette (blue) and the consensus sequence for *E. coli* BPD ABC transporters (red) (Boos & Lucht, 1996) of the 200 aa domain containing the ABC. The ATP-binding motif domains (A) and (B) are underlined, and the linker peptide (L) is identified. The proteins aligned are from the following organisms: GepD (*D. nodosus*), HI0354 (*H. influenzae*), TauB (*E. coli*), AtsC (*P. aeruginosa*), AF0092 (*A. fulgidus*), YgaL (*B. subtilis*), ShiC (*S. hydrodystenterie*). Amino acids that are identical in the six sequences aligned with *gepD* (though in some cases different from the consensus) are indicated by asterisks; conserved amino acid residues are indicated by a period; 10 aa intervals are shown (v).

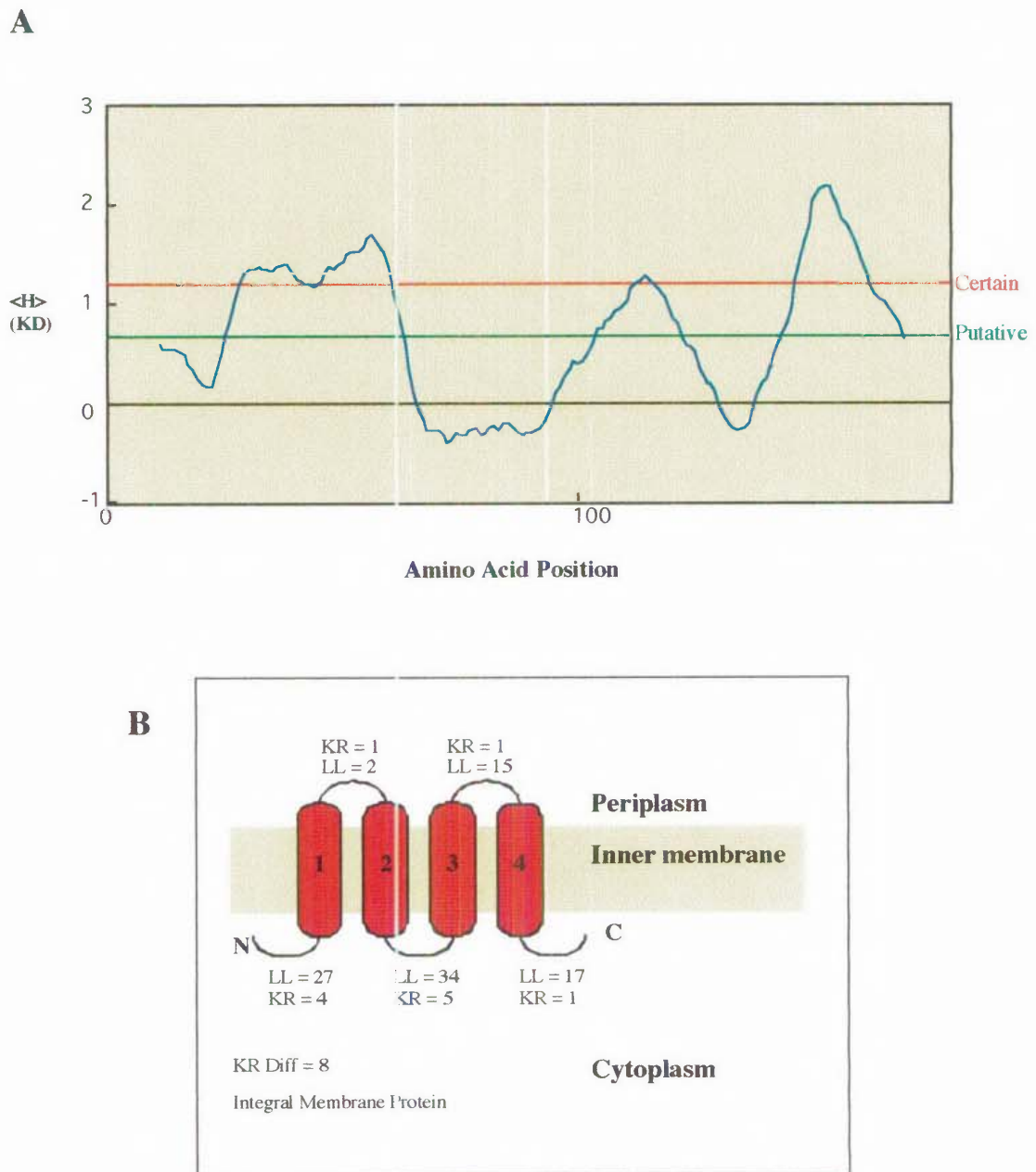


Figure 5.7: Hydrophobicity ($\langle H \rangle$) plot (A) and predicted topology (B) of the GepE putative BPD ABC transporter membrane-associated protein using the computer program TopPred II (Claros, 1996). Hydrophilicity is calculated according to the algorithm of Kyte-Doolittle (KD) (Kyte & Doolittle, 1982). Transmembrane domains are numbered 1 to 4. Other features are as follows: LL indicates the number of amino acids in the loop; KR indicates the number of lysine and arginine residues respectively; KR Diff indicates the positive charge difference; N- and C- indicate the amino and carboxy termini of the protein. Cytoplasmic and periplasmic faces and the inner membrane are also indicated.

(Fath & Kolter, 1993; Higgins, 1992). Alternatively, the observation that GepE does not contain the expected number of transmembrane subunits may also be a limitation of the program used to predict the topology of the GepE protein.

Despite the high degree of similarity between the ATP-binding cassette motifs of different ABC transporter proteins in general, there is often little sequence similarity between the membrane-associated proteins of different ABC transporters. It has been proposed that the folding of a membrane protein is likely to be less sensitive to amino acid substitutions than a soluble protein, and hence hydrophobic domains tend to be more divergent than hydrophilic ones (Higgins, 1992).

A ClustalW alignment of the putative ABC transporter membrane-associated protein, GepE, with similar transmembrane proteins, shows that although all have similarity to GepE (Table 5.2), they are divergent from each other. Despite this divergence, there are seven amino acid residues that are absolutely conserved including three glycine, three proline and one phenylalanine (Figure 5.8). This suggests that these particular residues may have an important functional or structural role in these related proteins.

Glycine, is often found at turns in the protein secondary structure, primarily because where it is present in a protein the minimal steric hindrance of the glycine side chain results in more structural flexibility than other amino acids. Proline is the opposite of glycine, since the secondary amino group is held in a rigid conformation that decreases the flexibility of the protein at the point at which it is present (Lehninger, Nelson & Cox, 1993). The prediction of α -helices, β -sheets, and turns in GepE was done using the GCG peptidestructure program (Jameson & Wolf, 1988) (Figure 5.8) available through ANGIS. The position of conserved proline and glycine residues in GepE relative to potential

secondary structures is consistent with the hypothesis that these residues may be important to the structural and therefore functional integrity of the protein.

In the membrane-associated protein of the BPD transport systems there is a single conserved sequence between 94 and 115 aa from the C-terminus feature which has the consensus sequence EAA---G-----I-LP (Dassa & Hofnung, 1985), called the EAA-loop. More recent studies of forty-seven membrane subunits led to the identification of a similar sequence in all proteins but with variations in some positions of this consensus sequence (Benner *et al.*, 1994). Similar variations are evident in the ClustalW alignment of GepE and related sequences (Figure 5.8). The exact function of the EAA-loop has not yet been elucidated, however it has been proposed that since it is located on the cytoplasmic side of the membrane it may interact with a ligand or protein, possibly the ATP-binding subunit of the BPD ABC transporter (Kerppola *et al.*, 1991).

5.2.4 A role for *gepC*?

In two-thirds of BPD transporters so far characterised, there are two separate membrane-associated proteins (Boos & Lucht, 1996), and it is assumed that these two membrane-associated proteins form a heterodimer within the transport complex. In the other third, the membrane-associated protein is thought to form a homodimer (Boos & Lucht, 1996; Fath & Kolter, 1993; Higgins, 1992).

In the *gep* operon only GepE had similarity to membrane-associated proteins of BPD importer proteins. However, computer aided analysis of the hydrophobicity and predicted topology of the GepC protein indicates that the GepC protein is also an integral membrane protein (Figure 5.9). The GepC has three potential transmembrane domains, that are separated by two hydrophilic loops, and an N- terminus which is located on the cytoplasmic face, whilst the C-terminus is found on the periplasmic side of the membrane.